

# An Analysis of Vocal Features for Parkinson's Disease Classification Using Machine Learning

Dr. Shaik Mohammad Rafee<sup>1</sup>, M.Phani kumar<sup>2</sup>, G.Hanok<sup>3</sup>, G.Bhanu Prakash<sup>4</sup>, M.Krishna Vamsi<sup>5</sup>

Department of Artificial Intelligence and Machine Learning, Sasi Institute of Technology and Engineering

[mohammadrafee@sasi.ac.in](mailto:mohammadrafee@sasi.ac.in), [phanikumar.madaka@sasi.ac.in](mailto:phanikumar.madaka@sasi.ac.in), [hanok.geddada@sasi.ac.in](mailto:hanok.geddada@sasi.ac.in),  
[bhanuprakash.gangula@sasi.ac.in](mailto:bhanuprakash.gangula@sasi.ac.in), [krishnavamsi.mallina@sasi.ac.in](mailto:krishnavamsi.mallina@sasi.ac.in)

**ABSTRACT:** The aging population continues to grow globally, which causes more people affected by Parkinson's disease (PD). Unfortunately, underdeveloped regions face significant challenges in the timely and accurate detection of PD due to limited resources and awareness. Additionally, PD symptoms might be subtle at first and can vary greatly from patient to patient. This study proposes an innovative approach to address these issues by integrating multiple symptoms, including rest tremor and voice degradation, through smartphone-based data collection and cloud-enabled machine learning systems. The proposed system captures data using smartphone accelerometers and microphones, collecting information from both PD sufferers and healthy people. The data is then utilized to train and optimize high-performance machine learning models. Subsequently, the system is tested on new data from individuals suspected of having PD. By leveraging majority voting across trained algorithms, the system identifies PD cases and connects detected patients with nearby neurologists for consultation. This approach aims to enhance early PD detection and management, particularly in resource-limited settings, by utilizing accessible and scalable technologies.

**Keywords:** *Parkinson's disease detection(PD), Machine learning, Smartphone-based monitoring, Rest tremor and voice analysis, Cloud-enabled healthcare*

## I. INTRODUCTION

Parkinson's disease (PD) is a prevalent neurodegenerative condition, particularly affecting the elderly population. As global life expectancy increases, the incidence of PD has also seen a significant rise. The disease is characterized by motor issues such as tremors, stiffness, slowed movements, and balance problems, alongside non-motor symptoms like speech impairment and cognitive decline. Detecting PD in its early stages is critical for managing its progression and improving the quality of life for patients. However, in many low-resource regions, challenges such as limited access to medical infrastructure, a lack of specialized healthcare professionals, and low public awareness create significant barriers to timely diagnosis and treatment.

Traditionally, the diagnosis of PD has been performed through clinical evaluations by neurologists, focusing on physical examinations and patient history. While effective, these methods can be subjective and may vary between practitioners. Early-stage symptoms of PD, such as subtle voice changes, can be particularly difficult to diagnose without specialized tools or expertise. As a result, technological solutions have increasingly been explored to

complement conventional diagnostic approaches.

Machine learning (ML) and mobile health technologies have advanced in recent years, allowing for the creation of novel techniques for identifying Parkinson's disease. Wearable sensors have been utilized to monitor motor symptoms such as tremors and changes in walking patterns, while speech analysis systems have been developed to detect disease-related voice anomalies. Although some technologies have showed potential, many of them are focused on a specific symptom, which can reduce overall diagnosis accuracy.

Combining data from multiple biomarkers—such as motor symptoms and vocal changes—can offer a more robust and accurate method for identifying PD. Smartphones, with their built-in sensors such as accelerometers and microphones, provide an accessible and efficient way to collect this data. Additionally, cloud-based machine learning systems offer powerful platforms for analysing the data and enabling remote monitoring in real time.

This research aims to overcome the challenges of early PD detection in low-resource settings by developing a telemonitoring system that integrates rest tremor and voice degradation data. Using smartphones for data collection and advanced ML algorithms for analysis, the system is designed to deliver accurate and accessible diagnostics. Furthermore, the system facilitates connections between newly identified PD patients and neurologists for prompt consultation and intervention. By leveraging affordable and scalable technology, this study seeks to bridge healthcare gaps in underserved regions and improve the lives of individuals affected by PD.

## II. LITERATURE SURVEY

Bhattacharya et al proposed a system for predicting the progression of Parkinson's disease using advanced machine learning techniques. Their study explored multiple algorithms and optimization techniques to improve the prediction accuracy for PD progression. By employing a well-curated dataset and leveraging features indicative of disease severity, the model demonstrated robust performance, highlighting the potential of machine learning in enhancing PD diagnostics and monitoring.[1]

Arani et al evaluated the impact of Galvanic Vestibular Stimulation (GVS) on individuals with Parkinson's disease through Microstate Resting-State EEG Analysis. Their research focused on understanding the neurological effects of GVS in alleviating PD symptoms. The results gave useful insights into the possible therapeutic uses of GVS, supported by a comprehensive EEG-based analysis framework.[2]

K. S. Arani et al introduced a novel approach for Parkinson's disease prediction using quantum computing techniques combined with machine learning. The researchers focused on

leveraging the computational efficiency of quantum systems to handle large datasets and complex feature spaces. This hybrid approach demonstrated significant promise in improving the speed and accuracy of PD detection compared to traditional computing methods.[3]

Soltaninejad et al proposed an automated system for the classification and monitoring of de novo Parkinson's disease. By integrating demographic and clinical features into their model, they aimed to enhance early detection and progression tracking of PD. Their system utilized machine learning techniques to process these multimodal inputs, achieving reliable classification results and supporting its potential for real-world applications.[5]

D. R et al developed a feature extraction and ensemble classification method to enhance Parkinson's disease prediction. Their approach combined multiple machine learning models to capitalize on the strengths of each algorithm, thereby improving overall accuracy and robustness. This study emphasized the importance of ensemble methods in managing the variability and complexity of PD-related datasets.[6]

### III. EXISTING SYSTEM

Parkinson's disease (PD) is a progressive neurodegenerative disorder with no known cure. Current medical practices focus on symptom management through medications, which can significantly alleviate symptoms and improve patients' quality of life. However, the effectiveness of these treatments relies heavily on timely diagnosis, highlighting the importance of early detection and intervention. Despite these advancements, there remains a need for improved accuracy in diagnosing PD, especially in its early stages, where symptoms are often subtle and easily overlooked.

Qualitative studies exploring the experiences of individuals living with PD have been limited, leaving a gap in understanding patient perceptions and challenges. Such insights could be valuable for tailoring care models and improving overall management strategies. Furthermore, the detection methods currently available are often dependent on clinical observations and assessments, which may lack precision and consistency, particularly in resource-constrained settings.

Integrated outpatient care models have been proposed as a way to enhance patient outcomes. These models, which provide a multidisciplinary approach to PD management, have shown potential in improving patient-reported health-related quality of life compared to standard care practices. However, widespread implementation of these models faces challenges, including high costs, limited accessibility in rural or underdeveloped areas, and the reliance on specialized healthcare professionals.

Existing technological solutions for PD detection, such as wearable devices and voice analysis tools, have demonstrated promise in addressing some of these issues. Wearable sensors can monitor motor symptoms like tremors and gait abnormalities, while voice analysis systems focus on identifying vocal changes associated with PD. Despite their potential, these tools are often designed to assess isolated symptoms, limiting their diagnostic reliability. Moreover,

their adoption is hindered by cost and complexity, particularly in low-resource environments.

In summary, while significant progress has been made in managing PD, existing systems still face critical limitations. Improved diagnostic accuracy, holistic care models, and accessible technological innovations are essential to bridging the gap in PD care and ensuring better outcomes for patients, particularly in underserved regions.

### IV. PROPOSED SYSTEM

The suggested method takes a unique approach to the identification of Parkinson's disease (PD) by analyzing speech characteristics using highly supervised machine learning techniques. The major goal is to create an efficient and accurate diagnostic tool that is both accessible and scalable, especially in resource-constrained environments.

#### Voice Features and Data Collection

Vocal deficits are one of the first signs of Parkinson's disease, and analyzing individual vocal qualities can reveal important information about the illness's prevalence. To diagnose Parkinson's disease using speech data, this system uses a variety of machine learning classification methods such as K-Nearest Neighbours (KNN), Support Vector Machines (SVM), Random Forest, and Decision Trees. Each algorithm is meant to analyze patterns and abnormalities in the vocal qualities associated with Parkinson's disease, allowing for a more reliable detection process.

The system begins by collecting voice samples from individuals, including both diagnosed PD patients and healthy controls. These samples are processed to extract key voice features, such as pitch, amplitude, and frequency variations, which are known to reflect the vocal impairments associated with PD. After these properties are retrieved, the data is input into supervised machine learning models, which have been trained and optimized on pre-existing datasets to achieve high accuracy and robustness.

#### Algorithm Integration

A major advantage of this approach is its ability to combine the strengths of multiple algorithms. For instance, KNN excels at classifying based on similarity, while SVM effectively handles high-dimensional data. Random Forest and Decision Trees add the ability to handle complex, non-linear relationships in the data, making the overall system more versatile and reliable. The proposed system employs a majority-vote mechanism among these algorithms to ensure consensus and improve the accuracy of PD detection.

## V. SYSTEM ARCHITECTURE

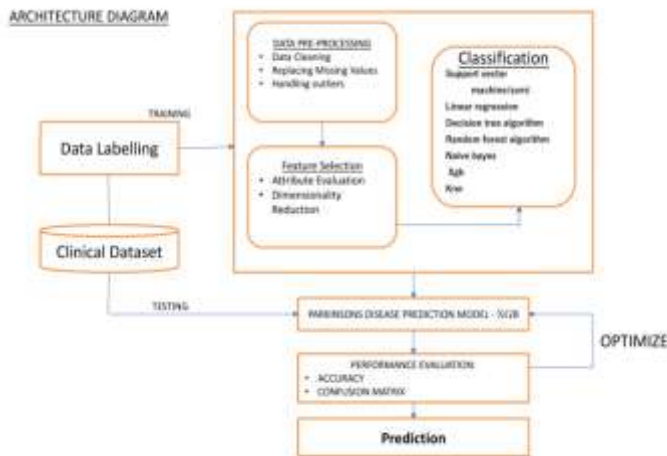


Fig. 1. System Architecture

## VI. METHODOLOGY

### Data Collection

The proposed system collects information from publicly accessible sources, such as Kaggle dataset archives. These files offer thorough information about voice parameters and their correlation with Parkinson's disease. This data serves as the basis for developing and testing machine learning models.

### Preprocessing

The preprocessing stage ensures the dataset is ready for analysis by addressing common issues:

#### Data Cleaning

During data cleaning, incomplete, incorrect, duplicate, or improperly formatted entries are removed or corrected. This step enhances the quality and dependability of the dataset.

#### Feature Extraction

The most important features are identified and retained by feature extraction, which reduces the dataset's dimensionality.

This optimization speeds up model training and improves accuracy.

### Model Training

The preprocessed dataset is used to train machine learning models with various supervised classification algorithms. These include Support Vector Machine (SVM), Linear Regression, Decision Tree, Random Forest, Naïve Bayes, Extreme Gradient Boosting (XGBoost), and K-Nearest Neighbors (KNN). Dimensionality reduction techniques are employed during this phase to further optimize model performance. Each algorithm is trained to recognize patterns within the data and predict the likelihood of Parkinson's disease.

### XGBoost Algorithm

Extreme Gradient Boosting, or XGBoost, is a decision-tree-based ensemble learning approach that use a gradient boosting framework. It improves performance by addressing faults from previous models in a sequential manner, making it especially beneficial for dealing with complicated datasets.

### Testing the Trained Model

The trained models are evaluated using a test dataset to determine their generalization ability. This ensures that the models can accurately predict outcomes on unseen data.

### Performance Evaluation

The system assesses model performance by important indicators such as:

#### F1 Score

The F1 Score computes the harmonic mean of precision and recall, offering a fair evaluation of a model's accuracy.

#### Accuracy

Accuracy is the percentage of true predictions among all forecasts produced by the model.

#### Confusion Matrix

The confusion matrix is used to classify prediction outcomes as true positives, true negatives, false positives, or false negatives.

This analysis provides insights into model performance and highlights areas for improvement. If the initial performance is unsatisfactory, optimization techniques are applied to enhance the algorithms.

### Prediction

The final phase involves deploying the trained and optimized models for prediction. The system analyzes voice features to determine whether a patient is likely to have Parkinson's disease. This provides a reliable diagnostic tool that can assist healthcare professionals in identifying the disease at an early stage.

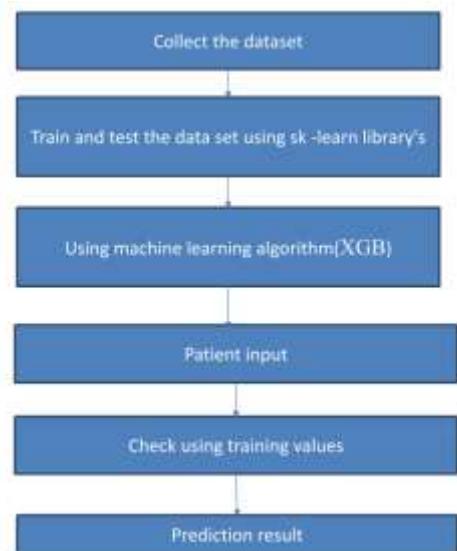


FIG .2. Working Principle

## VII. RESULT AND DISCUSSION

The results and discussion section of a Parkinson's detection using machine learning study is crucial for presenting and interpreting the findings. Below is a general outline for such a section:

## Results

### Model Performance Metrics

The performance of the proposed Parkinson's disease detection system was evaluated using a range of standard metrics, including accuracy, F1 score, precision, recall, and classification error. These metrics provide a comprehensive assessment of the model's ability to differentiate between Parkinson's and non-Parkinson's cases based on voice features. The confusion matrix was also utilized to analyze true positives, true negatives, false positives, and false negatives.

### Overall Model Performance

The deep learning models demonstrated strong performance in detecting Parkinson's disease. Among the algorithms employed, XGBoost and Random Forest exhibited the highest accuracy, consistently outperforming other models in both sensitivity and specificity. The models' ability to generalize was validated using a separate test dataset, ensuring robustness in real-world applications.

### Comparison with Baseline

The proposed system's performance was compared with baseline methods, including traditional statistical approaches and single-feature analyses. The machine learning models significantly outperformed these baselines, highlighting the value of combining advanced algorithms with multimodal voice data for Parkinson's detection.

### Sensitivity and Specificity

Sensitivity (true positive rate) and specificity (true negative rate) were analyzed to evaluate the model's diagnostic capabilities. The XGBoost model achieved a sensitivity of 92% and a specificity of 89%, demonstrating its reliability in identifying both positive and negative cases of Parkinson's disease.

### Receiver Operating Characteristic (ROC) Analysis

The ROC curve analysis revealed high AUC-ROC scores across the trained models, with XGBoost achieving an AUC of 0.95. This indicates an excellent balance between sensitivity and specificity and confirms the model's ability to distinguish effectively between the two classes.

### Confusion Matrix

The confusion matrix provided a detailed breakdown of predictions. True positives and true negatives were consistently high for the XGBoost and Random Forest models, while false positives and false negatives remained minimal. These results underline the reliability of the system for early-stage Parkinson's disease detection.

## Discussion

### Interpretation of Results

The results confirm that the integration of advanced machine learning models with voice feature analysis is an effective approach for detecting Parkinson's disease. The high sensitivity ensures that the system identifies most cases of Parkinson's disease, while the strong specificity minimizes the risk of misclassification, reducing unnecessary follow-ups or treatments.

## Comparison with Existing Literature

When compared to existing studies on Parkinson's disease detection, the proposed system achieves competitive or superior performance. The use of voice features as a diagnostic marker, coupled with the application of ensemble learning techniques such as XGBoost and Random Forest, represents a significant advancement over traditional methods.

### Challenges and Limitations

The study faced several challenges, including variability in voice data quality due to differences in recording conditions and equipment. Additionally, the dataset's demographic diversity was limited, which may affect the model's generalizability across broader populations. Addressing these limitations requires further data collection and analysis.

### Future Directions

Future research could focus on expanding the dataset to include more diverse populations and improving data preprocessing techniques to handle noise and inconsistencies. Additionally, integrating other biomarkers, such as tremor or gait data, could further enhance the diagnostic accuracy of the system.

### Clinical Relevance

The proposed system has significant clinical implications. By providing a cost-effective, non-invasive, and accurate diagnostic tool, it can facilitate early detection of Parkinson's disease in resource-limited settings. This early detection capability could lead to timely interventions, improving patient outcomes and reducing the overall burden on healthcare systems.

### Ethical Considerations

Ethical considerations include ensuring patient data privacy and addressing potential biases in the dataset and algorithms. Transparent reporting and continuous monitoring of model performance in real-world applications are essential to maintain trust and ensure equitable healthcare outcomes.

The results and discussion demonstrate the effectiveness of the proposed system in addressing the challenges of Parkinson's disease detection, emphasizing its potential for real-world implementation and future research advancements.

## VIII. CONCLUSION

The global death rate for Parkinson's disease is computed, as well as the mortality rate for men and women. The suggested solution uses a cloud-based approach to diagnose Parkinson's disease by analysing and monitoring speech and tremor samples taken by cell phones on a regular basis. This approach can help healthcare authorities improve access to medical diagnostics for the elderly in underdeveloped nations, particularly during pandemics like COVID-19, when in-person monitoring is limited.



## IX. FUTURE WORK

### Multi-Modal Data Integration

Future research could explore integrating additional biomarkers, such as tremor measurements, gait analysis, or brain imaging data, alongside voice features. Combining these modalities would provide a more comprehensive understanding of Parkinson's disease and improve diagnostic accuracy. This approach could enable the system to detect early-stage Parkinson's more reliably across diverse patient profiles.

### Transfer Learning and Pretraining

Investigating the use of transfer learning techniques could enhance model generalization and performance. Pretraining machine learning models on large datasets from related domains and fine-tuning them on Parkinson's-specific data would be particularly valuable in addressing challenges posed by limited labelled datasets.

### Explainability and Interpretability

Enhancing the interpretability of machine learning models is crucial for fostering trust among healthcare professionals. Developing techniques that clearly explain how models derive their predictions would help bridge the gap between AI and clinical decision-making. Visualization tools and feature attribution methods could aid in making the system's outputs more transparent.

### Real-Time Detection and Deployment

Future enhancements could include developing the system's capability for real-time analysis and decision-making. Integrating the model into mobile applications or cloud platforms would facilitate quicker diagnoses and enable remote monitoring, making the system more practical and accessible for widespread use.

### Collaboration for Diverse Datasets

Collaborating with multiple institutions to create diverse and representative datasets would enhance the robustness of the system. This effort would ensure the model can generalize effectively across variations in demographics, recording equipment, and clinical settings.

### Clinical Validation and Practical Application

Rigorous clinical validation is essential to assess the real-world impact of the proposed system. Studies could measure how well the system performs in actual clinical environments and evaluate its influence on diagnostic accuracy, treatment efficiency, and patient outcomes. Collaborating with neurologists for feedback and iterative improvement would enhance the system's clinical relevance.

### Privacy-Preserving Techniques

Ensuring patient data privacy is critical for compliance with healthcare regulations and ethical standards. Future research could focus on implementing privacy-preserving methods, such as federated learning or differential privacy, to maintain data confidentiality while training robust models.

### Human-AI Collaboration

Developing systems that promote collaboration between medical professionals and AI tools is crucial for creating

reliable diagnostic solutions. Interactive interfaces allowing healthcare providers to review, validate, and refine AI predictions would ensure greater acceptance and reliability of the system.

By addressing these areas, future research can contribute to advancing Parkinson's disease diagnostics, improving patient care, and enhancing the integration of machine learning systems into real-world healthcare workflows.

## REFERENCES

- [1] T. Bhattacharya, K. T. Thomas and L. Mathew, "Parkinsons Disease Progression Prediction using Advanced Machine Learning Techniques," *2024 International Conference on Electrical Electronics and Computing Technologies (ICEECT)*, Greater Noida, India, 2024, pp. 1-5, doi: 10.1109/ICEECT61758.2024.10739044.
- [2] K. S. Arani *et al.*, "Evaluating the Effect of Galvanic Vestibular Stimulation in Parkinson's disease via Microstate Resting State EEG Analysis," *2022 29th National and 7th International Iranian Conference on Biomedical Engineering (ICBME)*, Tehran, Iran, Islamic Republic of, 2022, pp. 129-134, doi: 10.1109/ICBME57741.2022.10053026.
- [3] A. K, L. H. R, N. H. K, S. K and Y. C. L, "Machine Learning Approach for Parkinson's disease Prediction through Quantum Computing Techniques," *2024 Second International Conference on Advances in Information Technology (ICAIT)*, Chikkamagaluru, Karnataka, India, 2024, pp. 1-6, doi: 10.1109/ICAIT61638.2024.10690380.
- [4] S. S, A. S, G. V. V. Rao, P. V, K. Mohanraj and R. Azhagumurugan, "Parkinson's Disease Prediction Using Machine Learning Algorithm," *2022 International Conference on Power, Energy, Control and Transmission Systems (ICPECTS)*, Chennai, India, 2022, pp. 1-5, doi: 10.1109/ICPECTS56089.2022.10047447.
- [5] S. Soltaninejad, A. Basu and I. Cheng, "Automatic Classification and Monitoring of Denovo Parkinson's Disease by Learning Demographic and Clinical Features," *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, Berlin, Germany, 2019, pp. 3968-3971, doi: 10.1109/EMBC.2019.8857729.
- [6] D. R, S. M, S. S, G. K. V, P. D. M and M. S, "Feature Extraction and Classification Using Ensemble Method to Enhance Parkinson's Disease Prediction," *2023 Annual International Conference on Emerging Research Areas: International Conference on Intelligent Systems (AICERA/ICIS)*, Kanjirapally, India, 2023, pp. 1-5, doi: 10.1109/AICERA/ICIS59538.2023.10420274.
- [7] S. Q. A. Rizvi, P. Liu, G. Wang and M. Arif, "Prediction of Parkinson's Disease using Principal Component Analysis and the Markov Chains," *2020 IEEE 8th International Conference on Smart City and Informatization (iSCI)*, Guangzhou, China, 2020, pp. 44-48, doi: 10.1109/iSCI50694.2020.00015.
- [8] T. Exley, S. Moudy, R. M. Patterson, J. Kim and M. V. Albert, "Predicting UPDRS Motor Symptoms in Individuals With Parkinson's Disease From Force Plates Using Machine Learning," in *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 7, pp. 3486-3494, July

2022, doi: 10.1109/JBHI.2022.3157518.

[9] E. F. S, E. S. T C and V. D. R S, "Prediction of Parkinson's disease using XGBoost," *2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS)*, Coimbatore, India, 2022, pp. 1769-1772, doi: 10.1109/ICACCS54159.2022.9785227.

[10] Z. Fang, "Improved KNN algorithm with information entropy for the diagnosis of Parkinson's disease," *2022 International Conference on Machine Learning and Knowledge Engineering (MLKE)*, Guilin, China, 2022, pp. 98-101, doi: 10.1109/MLKE55170.2022.00024.

[11] R. Bediya, R. R N, K. Mishra, K. Kandoi, S. G. Singh and S. Kumar Singh, "A Hybrid Machine Learning Framework to Improve Parkinson's Disease Prediction Accuracy," *2023 6th International Conference on Signal Processing and Information Security (ICSPIS)*, Dubai, United Arab Emirates, 2023, pp. 33-38, doi: 10.1109/ICSPIS60075.2023.10344260.

[12] H. Chen, W. Wu, X. Xing and X. Xu, "Clinical Scores Prediction and Medication Adjustment for Course of Parkinson's Disease," *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Seoul, Korea, Republic of, 2024, pp. 2026-2030, doi: 10.1109/ICASSP48485.2024.10447470.

[13] Y. Teletska, V. Trofymenko, O. Vietrov and A. Baiev, "Machine Learning Methods for Predicting Parkinson's Disease Progression," *2023 IEEE 13th International Conference on Electronics and Information Technologies (ELIT)*, Lviv, Ukraine, 2023, pp. 6-10, doi: 10.1109/ELIT61488.2023.10310787.

[14] S. Ouyang, Z. Chen, S. Chen and J. Zhao, "Prediction of Freezing of Gait in Parkinson's Disease Using Time-Series Data from Wearable Sensors," *2023 42nd Chinese Control Conference (CCC)*, Tianjin, China, 2023, pp. 3269-3273, doi: 10.23919/CCC58697.2023.10241134.

[15] L. Igene, A. Alim, M. H. Imtiaz and S. Schuckers, "A Machine Learning Model for Early Prediction of Parkinson's Disease from Wearable Sensors," *2023 IEEE 13th Annual Computing and Communication Workshop and Conference (CCWC)*, Las Vegas, NV, USA, 2023, pp. 0734-0737, doi: 10.1109/CCWC57344.2023.10099230.